# AutoShuffleNet: Learning Permutation Matrices via an Exact Lipschitz Continuous Penalty in Deep Convolutional Neural Networks

Jiancheng Lyu, Shuai Zhang, Yingyong Qi Qualcomm AI Research jianlyu,shuazhan,yingyong@qti.qualcomm.com Jack Xin
University of California, Irvine
jack.xin@uci.edu

# **ABSTRACT**

ShuffleNet is a state-of-the-art light weight convolutional neural network architecture. Its basic operations include group, channelwise convolution and channel shuffling. However, channel shuffling is manually designed on empirical grounds. Mathematically, shuffling is a multiplication by a permutation matrix. In this paper, we propose to automate channel shuffling by learning permutation matrices in network training. We introduce an exact Lipschitz continuous non-convex penalty so that it can be incorporated in the stochastic gradient descent to approximate permutation at high precision. Exact permutations are obtained by simple rounding at the end of training and are used in inference. The resulting network, referred to as AutoShuffleNet, achieved improved classification accuracies on data from CIFAR-10, CIFAR-100 and ImageNet while preserving the inference costs of ShuffleNet. In addition, we found experimentally that the standard convex relaxation of permutation matrices into stochastic matrices leads to poor performance. We prove theoretically the exactness (error bounds) in recovering permutation matrices when our penalty function is zero (very small). We present examples of permutation optimization through graph matching and two-layer neural network models where the loss functions are calculated in closed analytical form. In the examples, convex relaxation failed to capture permutations whereas our penalty succeeded.

# **KEYWORDS**

ShuffleNet; Permutation; Lipschitz Continuous Penalty; Convolutional Neural Network

# **ACM Reference Format:**

Jiancheng Lyu, Shuai Zhang, Yingyong Qi and Jack Xin. 2020. AutoShuf-fleNet: Learning Permutation Matrices via an Exact Lipschitz Continuous Penalty in Deep Convolutional Neural Networks. In *Proceedings of the 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '20), August 23–27, 2020, Virtual Event, CA, USA*. ACM, New York, NY, USA, 9 pages. https://doi.org/10.1145/3394486.3403103

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

KDD '20, August 23–27, 2020, Virtual Event, CA, USA © 2020 Association for Computing Machinery. ACM ISBN 978-1-4503-7998-4/20/08...\$15.00 https://doi.org/10.1145/3394486.3403103

### 1 INTRODUCTION

Light convolutional deep neural networks (LCNN) are attractive in resource limited conditions for delivering high performance at low costs. Some of the state-of-the-art LCNNs in computer vision are ShuffleNet [14, 23], IGC (Interleaved Group Convolutions, [22]) and IGCV3 (Interleaved Low-Rank Group Convolutions,[17]). A noticeable feature in the design is the presence of permutations for channel shuffling in between separable convolutions. The permutations are hand-crafted by designers outside of network training however. A natural question is whether the permutations can be learned like network weights so that they are optimized based on training data. An immediate difficulty is that unlike weights, permutations are highly discrete and incompatible with the stochastic gradient descent (SGD) methodology that is continuous in nature. To overcome this challenge, we introduce an exact Lipschitz continuous non-convex penalty so that it can be incorporated in SGD to approximate permutation at high precision and low overhead. Consequently, exact permutations are obtained by simple rounding at the end of network training with negligible drop of classification accuracy. To be specific, we shall work with ShuffleNet architecture [14, 23]. Our approach extends readily to other LCNNs.

**Related Work.** Permutation optimization is a long standing problem arising in operations research, graph matching among other applications [3, 8]. Well-known examples are linear and quadratic assignment problems [18]. Graph matching is a special case of quadratic assignment problem, and can be formulated over  $N \times N$  permutation matrices  $\mathcal{P}^N$  as:

$$\min_{\pi \in \mathcal{P}^N} \|A - \pi B \pi^T\|_F^2,$$

where A and B are the adjacency matrices of graphs with N vertices, and  $\|\cdot\|_F$  is the Frobenius norm. A popular way to handle  $\mathcal{P}^N$  is to relax it to the Birkhoff polytope  $\mathcal{D}^N$ , the convex hull of  $\mathcal{P}^N$ , leading to a convex relaxation. The explicit realization of  $\mathcal{D}^N$  is the set of doubly stochastic matrices

$$\mathcal{D}^N = \{ M \in \mathbf{R}^{N \times N} : M\mathbf{1} = \mathbf{1}, M^T\mathbf{1} = \mathbf{1}, M \ge 0 \},$$

where  $\mathbf{1}=(1,1,,\cdots,1)^T\in \mathbf{R}^N$ . An approximate yet simpler way to treat  $\mathcal{D}^N$  is through the classical first order matrix scaling algorithm, e.g. the Sinkhorn method, see [16] and its recent applications [7, 15]. Though in principle such algorithm converges, the cost can be quite high when iterating many times, which causes a bottleneck effect [12]. A non-convex and more compact relaxation of  $\mathcal{P}^N$  is by a sorting network [12] which maps the box  $[0,1]^N$  into a manifold that sits inside  $\mathcal{D}^N$  and contains  $\mathcal{P}^N$ . Yet another method is path following algorithm [21] which seeks solutions under concave relaxations of  $\mathcal{P}^N$  by solving a linear interpolation

of convex-concave problems (starting from the convex relaxation). Permutation learning via continuous approximation has been studied in visual data recovery [4]. None of the existing relaxations are exact. In the context of improving ShuffleNet, HadaNet [24] uses Hadamard matrices (H) to define a class of structured convolution as the product  $H^T \times$  group convolution  $\times H$  and generalize shuffled group convolution of ShuffleNet. However, the inference cost of HadaNet is much higher than that of ShuffleNet, and relies on special hardware for speedup. Hadamard matrices are constructed to date for certain special orders such as powers of 2, and conjectured to exist for multiples of 4. In particular, they are not applicable to odd channel/group numbers.

Contribution. Our non-convex relaxation is a combination of matrix  $\ell_{1-2}$  penalty function and  $\mathcal{D}^N$ . The  $\ell_{1-2}$  (the difference of  $\ell_1$ and  $\ell_2$  norms) has been proposed and found effective in selecting sparse vectors under nearly degenerate linear constraints [6, 20]. The matrix version is simply a sum of  $\ell_{1-2}$  over all row and column vectors. We prove that the penalty is zero when applied to a matrix in  $\mathcal{D}^N$  if and only if the matrix is a permutation matrix. Thanks to the  $\mathcal{D}^N$  constraint, the penalty function is Lipschitz continuous (almost everywhere differentiable). This allows the penalty to be integrated directly into SGD for learning permutation in LCNNs. As shown in our experiments on CIFAR-10, CIFAR-100 and ImageNet data sets, the closeness to  $\mathcal{P}^N$  turns out to be remarkably small at the end of network training so that a simple rounding has negligible effect on the validation accuracy. We also found that convex relaxation by  $\mathcal{D}^N$  fails to capture good permutations for LCNNs. We observed experimentally that a random shuffle could perform better than manual shuffle, but the learned shuffle consistently achieved the best results. To our best knowledge, this is the first time permutations have been successfully learned for the architecture selection of deep CNNs to improve hand-crafted permutations. Moreover, our AutoShuffleNet preserves the inference cost of ShuffleNet for any channel/group numbers.

**Outline.** In section 2, we introduce exact permutation penalty, and prove its closeness to permutation matrices when the penalty values are small, as observed in the experiments. We also present the training algorithm combining thresholding and matrix scaling to approximate projection onto  $\mathcal{P}^N$  for SGD. In section 3, we analyze two permutation optimization problems to show the utility of our penalty. In a 2-layer neural network regression model with short cut (identity map), convex relaxation does not give the optimal permutation even with additional rounding while our penalty can. In section 4, we show experimental results on consistent improvement of auto-shuffle over hand-crafted shuffle on data from CIFAR-10, CIFAR-100 and ImageNet. Conclusion is in section 5.

# 2 PERMUTATION, MATRIX $\ell_{1-2}$ PENALTY AND EXACT RELAXATION

The channel shuffle operation in ShuffleNet [14, 23] can be represented as multiplying the feature map in the channel dimension by a permutation matrix M. The permutation matrix M is a square binary matrix with exactly one entry of one in each row and each column and zeros elsewhere. In the ShuffleNet architecture [14, 23], M is preset by the designers and will be called "manual". In this work, we propose to learn an automated permutation matrix M

through network training, hence removing the human factor in its selection towards a more optimized shuffle. Since permutation is discrete in nature and too costly to enumerate, we propose to approach it by adding a matrix generalization of the  $\ell_{1-2}$  penalty [6, 20] to the network loss function in the stochastic gradient descent based training.

Specifically for  $M = (m_{ij}) \in \mathbb{R}^{N \times N}$ , the proposed continuous matrix penalty function is

$$P(M) := \sum_{i=1}^{N} \left[ \sum_{j=1}^{N} \left| m_{ij} \right| - \left( \sum_{j=1}^{N} m_{ij}^{2} \right)^{1/2} \right] + \sum_{j=1}^{N} \left[ \sum_{i=1}^{N} \left| m_{ij} \right| - \left( \sum_{i=1}^{N} m_{ij}^{2} \right)^{1/2} \right], \tag{1}$$

in conjunction with the doubly stochastic constraint:

$$m_{ij} \ge 0, \ \forall (i,j); \ \sum_{i=1}^{N} m_{ij} = 1, \ \forall j; \ \sum_{j=1}^{N} m_{ij} = 1, \ \forall i.$$
 (2)

REMARK 1. When the constraints in (2) hold, 
$$\sum_{i=1}^{N} |m_{ij}|$$
 and  $\sum_{i=1}^{N} |m_{ij}|$ 

in P(M) can be removed. However, in actual computation, the two equality constraints of (2) only hold approximately, so the full expression in (1) is necessary.

Remark 2. Thanks to (2), we see that the penalty function P(M) is actually Lipschitz continuous in M as  $\sum_{j=1}^{N} m_{ij}^2 \neq 0$ ,  $\forall i$ , and  $\sum_{i=1}^{N} m_{ij}^2 \neq 0$ 

 $0, \forall j.$  Although there are alternative penalties, we choose P(M) because it is simple, effective, and integrated well with SGD.

THEOREM 1. A square matrix M is a permutation matrix if and only if P(M) = 0, and the doubly stochastic constraint (2) holds.

PROOF.  $(\Rightarrow)$  Since a permutation matrix consists of columns (rows) with exactly one entry of 1 and the rest being zeros, each term inside the outer sum of P(M) equals zero, and clearly (2) holds.  $(\Leftarrow)$  By the elementary inequality,

$$\left(\sum_{j=1}^{N} |m_{ij}|\right) - \left(\sum_{j=1}^{N} m_{ij}^{2}\right)^{1/2} \ge 0, \quad \forall i,$$

with equality if and only if the row-wise cardinalty is 1:

$$|\{j: m_{ij} \neq 0\}| = 1, \forall i.$$
 (3)

This is because the mixed product terms like  $2 |m_{ij} m_{ij'}|$   $(j \neq j')$  in  $(\sum_{j=1}^{N} |m_{ij}|)^2$  must be all zero to match  $\sum_{j=1}^{N} m_{ij}^2$ . It only happens when equation (3) is true. Likewise,

$$\sum_{i=1}^{N} |m_{ij}| - \left(\sum_{i=1}^{N} m_{ij}^{2}\right)^{1/2} \ge 0, \quad \forall j,$$

with equality if and only if

$$|\{i: m_{ij} \neq 0\}| = 1, \forall j.$$

In view of (2), M is a permutation matrix.

The non-negative constraint in (2) is maintained throughout SGD by thresholding  $m_{ij} \rightarrow \max(m_{ij}, 0)$ . The normalization conditions in (2) are implemented sequentially once in one SGD iteration. Hence they are not strictly enforced. In theory, if the column/row normalization (divide each column/row by its sum) repeats sufficiently many times, the resulting matrices converge to (2), known as the Sinkhorn process [16]. We did not find much benefit to iterate more than once in terms of enhancing validation accuracy since the error in matrix scaling can be compensated in network weight adjustment during SGD.

The multiplication by M can be embedded in the network as a  $1 \times 1$  convolution layer with M initialized as absolute value of a random Gaussian matrix. After each weight update, we threshold the weights to  $[0, \infty)$ , normalize rows to unit lengths, then repeat on columns. Let *L* be the network loss function. The training minimizes the objective function:

$$f = f(w, M) := L(w) + \lambda \sum_{i=1}^{J} P(M_j),$$
 (4)

where J is the total number of "channel shuffle",  $M_i$ 's abbreviated as M, w is the network weight,  $\lambda$  a positive parameter. The training algorithm is summarized in Alg. 1. Introducing those  $1 \times 1$  convolutions and the penalty term results in little extra computation, so the training time is similar to training ShuffleNet. The  $\ell_1$  term in the penalty function P has standard sub-gradient, and the  $\ell_2$  term is differentiable away from zero, which is maintained in the Alg. 1 by SGD and normalization in columns and rows.  $\lambda$  is chosen to be  $10^{-3}$  or  $2 \times 10^{-3}$  so as to balance the contributions of the two terms

in (4) and drive 
$$\sum_{j=1}^{J} P(M_j)$$
 close to 0.

We shall see that the penalty P indeed gradually gets smaller during training (Fig. 8). Here we show a theoretical bound on the distance to  $\mathcal{P}^N$  when P is small and (2) holds approximately.

Theorem 2. Let the dimension N of a non-negative square matrix M be fixed. If  $P(M) = O(\epsilon)$ ,  $\epsilon \ll 1$ , and the doubly stochastic constraints are satisfied to  $O(\epsilon)$ , then there exists a permutation matrix  $P^*$  such that  $||M - P^*||_F = O(\epsilon)$ .

PROOF. It follows from  $P(M) = O(\epsilon)$  that

$$\left(\sum_{j=1}^{N}\left|\,m_{ij}\,\right|\right)-\left(\sum_{j=1}^{N}m_{ij}^{2}\right)^{1/2}=O(\epsilon),\ \, \forall i,$$

implying that:

$$|m_{ij} m_{ij'}| = O(\epsilon), \forall j \neq j', \forall i.$$
 (5)

On the other hand for  $\forall i$ :

$$\sum_{j=1}^{N} m_{ij} = 1 + O(\epsilon). \tag{6}$$

Let  $j^* = \operatorname{argmax}_{i} |m_{ij}|$ , at any i. It follows from (6) that

$$|m_{i,i^*}| \geq 1/N + O(\epsilon),$$

and from (5) that

$$m_{i,i'} = O(\epsilon), \quad \forall j' \neq j^*.$$

#### Algorithm 1 AutoShuffle Learning.

mini-batch loss function  $f_t(w, M)$ , t being the iteration index; learning rate  $\eta^t$  for (w, M); penalty parameter  $\lambda$  for P;

total iteration number Tn.

#### Start:

w: sample from unit Gaussian distribution;

M: sample from unit Gaussian distribution then take absolute value.

### WHILE t < Tn, DO:

- (1) Evaluate the mini-batch gradient  $(\nabla_w f_t, \nabla_M f_t)$  at  $(w^t, M^t)$ ; (2)  $w^{t+1} = w^t \eta^t \nabla_w f_t(w^t, M^t)$ ; // gradient update for
- (3)  $M^{t+1} = M^t \eta^t \nabla_M f_t(w^t, M^t)$ ; // gradient update for M
- (4)  $M^{t+1} \leftarrow \max(M^{t+1}, 0)$ ; // thresholding to enforce nonnegativity constraint
- (5) normalize each column of  $M^{t+1}$  by dividing the sum of entries in the column;
- (6) normalize each row of  $M^{t+1}$  by dividing the sum of entries in the row.

# END WHILE

**Output**:  $w^{Tn}$ ,  $M^{Tn}$ ; project each matrix  $M_j^{Tn}$  inside  $M^{Tn}$  to the nearest permutation matrix.

Hence each row of M is  $O(\epsilon)$  close to a unit coordinate vector, with one entry near 1 and the rest near 0. Similarly from

$$\sum_{i=1}^{N} \left| m_{ij} \right| - \left( \sum_{i=1}^{N} m_{ij}^2 \right)^{1/2} = O(\epsilon), \quad \forall j,$$

and 
$$\sum_{i=1}^{N} m_{ij} = 1 + O(\epsilon)$$
, we deduce that each column of  $M$  is  $O(\epsilon)$ 

close to a unit coordinate vector, with one entry near 1 and the rest near 0. Combining the two pieces of information above, we conclude that M is  $O(\epsilon)$  close to a permutation matrix.

The learned non-negative matrix *M* will be called a *relaxed shuffle* and rounded to the nearest permutation matrix to produce a final auto shuffle. Relaxed shuffle usually has better performance before rounding but the auto shuffle is desirable since it preserves the shuffle structure of the original ShuffleNet without incurring extra computation. Strictly speaking, this "rounding" involves finding the orthogonal projection to the set of permutation matrices, a problem called the linear assignment problem (LAP), see [1] and references therein. The LAP can be formulated as a linear program over the doubly stochastic matrices or constraints (2), and is solvable in polynomial time [1]. As we shall see later in Table 5, the relaxed shuffle comes amazingly close to an exact permutation in network learning. It is unnecessary to solve LAP exactly, indeed a simple rounding will do. AutoShuffleNet units adapted from ShuffleNet v1 [23] and ShuffleNet v2 [14] are illustrated in Figs. 1-2.

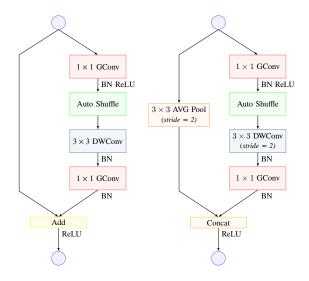


Figure 1: AutoShuffleNet units based on ShuffleNet v1.

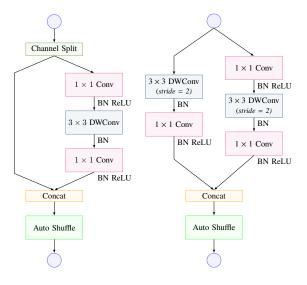


Figure 2: AutoShuffleNet units based on ShuffleNet v2.

# 3 PERMUTATION PROBLEMS UNSOLVABLE BY CONVEX RELAXATION

The doubly stochastic matrix condition (2) is a popular convex relaxation of permutation. However, it is not powerful enough to enable auto-shuffle learning as we shall see later. In this section, we present examples from permutation optimization to show the limitation of convex relaxation (2), and how our proposed penalty (1) can strengthen (2) to retrieve permutation matrices.

Let us recall the graph matching (GM) problem, see [1, 12, 13, 18, 21] and references therein. The goal is to align the vertices of two graphs to minimize the number of edge disagreements. Given a pair of n-vertex graphs  $G_A$  and  $G_B$ , with respective adjacency  $n \times n$  matrices A and B, the GM problem is to find a permutation matrix

Q to minimize  $\|AQ-QB\|_F^2.$  Let  $\Pi$  be the set of all permutation matrices, solve

$$Q^* := \operatorname{argmin}_{Q \in \Pi} \|AQ - QB\|_F^2.$$
 (7)

By algebraic identity

$$||AQ - QB||_F^2$$
= trace{ $(AQ - QB)^T (AQ - QB)$ }  
= trace( $A^T A$ ) + trace( $B^T B$ ) - 2trace( $A O B^T O^T$ ).

the GM problem (7) is same as

$$Q^* = \operatorname{argmin}_{Q \in \Pi} \operatorname{trace}((-A) Q B^T Q^T),$$

a quadratic assignment problem (QAP). The general QAP for two real square matrices A and B is [12, 18]:

$$Q^* = \operatorname{argmin}_{Q \in \Pi} \operatorname{trace}(A Q B^T Q^T).$$

The convex relaxed GM is:

$$Q_* := \operatorname{argmin}_{Q \in D^N} \|A \, Q - Q \, B\|_F^2.$$

As an instance of general QAP, let us consider problem (7) in case n = 2 for two real matrices:

$$A = \left[ \begin{array}{cc} a & b \\ c & d \end{array} \right], \quad B = \left[ \begin{array}{cc} a' & b' \\ c' & d' \end{array} \right].$$

If  $Q \in \mathcal{D}^2$ , then:

$$Q = \left[ \begin{array}{cc} q & 1-q \\ 1-q & q \end{array} \right], \quad q \in [0,1];$$

and

$$AQ - QB = \left[ \begin{array}{cc} (a-a')\,q + (b-c')\,q' & (b-b')\,q + (a-d')\,q' \\ (c-c')\,q + (d-a')\,q' & (d-d')\,q + (c-b')\,q' \end{array} \right].$$

where q' = 1 - q.

Example 1: Let

$$A = \left[ \begin{array}{cc} 1 & 2 \\ 3 & 1 \end{array} \right], \quad B = \left[ \begin{array}{cc} 0 & 2 \\ 3 & 1 \end{array} \right].$$

$$AQ - QB = \begin{bmatrix} 2q - 1 & 0 \\ 1 - q & 1 - q \end{bmatrix},$$

$$||AQ - QB||_F^2 = (2q - 1)^2 + 2(1 - q)^2 = 6q^2 - 8q + 3,$$

which is convex on [0, 1] and has minimum at  $q_* = 2/3$ . The convex relaxed matrix solution is:

$$Q_* = \left[ \begin{array}{cc} 2/3 & 1/3 \\ 1/3 & 2/3 \end{array} \right],$$

however, the permutation matrix solution  $Q^*$  is the 2 × 2 identity matrix at q = 1.

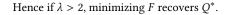
In the spirit of objective function (4), let us minimize

$$||AQ - QB||_E^2 + \lambda P(Q),$$

or equivalently minimize (after skipping additive constants in P)

$$F = F(q) := 6q^2 - 8q + 2 - 4\lambda(q^2 + (1-q)^2)^{1/2}$$
.

An illustration of F is shown in Fig. 3. The minimal point moves from the interior of the interval [0,1] when  $\lambda = 0.25$  (dashed line, top curve) to the end point 1 as  $\lambda$  increases to 1 (line-star, middle curve) and remains there as  $\lambda$  further increases to 2 (line-circle,



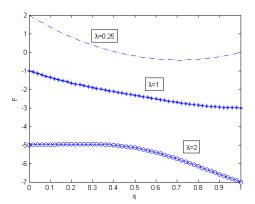


Figure 3: The function F(q) as penalty parameter  $\lambda$  varies from 0.25 (interior minimal point, dashed line, top) to 1 (linestar, middle) and 2 (line-circle, bottom). Minimal point occurs at q=1 in the latter two curves.

**Example 2:** Consider the adjacent matrix B(A) of an un-directed graph of 2 nodes and 1 edge (with a loop at node 1). An edge adds 1 and a loop adds 2 to an adjacent matrix.

$$A = \left[ \begin{array}{cc} 2 & 1 \\ 1 & 0 \end{array} \right], \quad B = \left[ \begin{array}{cc} 0 & 1 \\ 1 & 0 \end{array} \right].$$

Then:

$$AQ - QB = \begin{bmatrix} 2q & 2(1-q) \\ 0 & 0 \end{bmatrix},$$

$$||AQ - QB||_F^2 = 4[q^2 + (1-q)^2].$$

So

$$Q_* = Q(1/2) \neq Q^* = Q(0) = Q(1).$$

The P regularized objective function (modulo additive constants) is:

$$F = 4[q^2 + (1-q)^2] - 4\lambda(q^2 + (1-q)^2)^{1/2},$$

with  $F(0) = F(1) = 4 - 4\lambda$ . In view of

$$F'/4 = (2q-1)[2-\lambda/(q^2+(1-q)^2)^{1/2}],$$

two possible interior critical points are:

$$q = 1/2$$
 or  $q^2 + (1-q)^2 = \lambda^2/4$ . (8)

Since

$$\max_{q \in [0,1]} \{q^2 + (1-q)^2\} = 1,$$

the second equality in (8) is ruled out if  $\lambda > 2$ . Comparing

$$F(1/2) = 2 - 4\lambda 2^{-1/2} = 2(1 - \sqrt{2}\lambda)$$

with F(0), we see that the global minimal point does not occur at q = 1/2 if

$$1 - \sqrt{2}\lambda > 2 - 2\lambda$$
 or  $\lambda > 1/(2 - \sqrt{2}) \approx 1.7071$ 

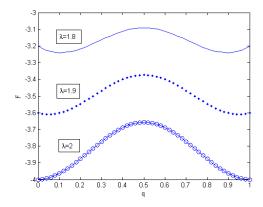


Figure 4: The function F(q) as penalty parameter  $\lambda$  varies from 1.8 (solid line, top) to 1.9 (dot, middle) and 2 (line-circle, bottom) where minimal points occur at q = 0, 1. Interior minimal points occur on [0, 1] when  $\lambda = 1.8, 1.9$ .

In Fig. 4, we show that two minimal points of F occur in the interior of (0,1) when  $\lambda = 1.8, 1.9$ , and transition to q = 0, 1, at  $\lambda = 2$ . When

$$\lambda^2/4 < \min_{q \in [0,1]} \{q^2 + (1-q)^2\} = 1/2$$

or  $\lambda < \sqrt{2}$ , the second equality in (8) cannot hold, *F* becomes convex with a unique minimal point at q = 1/2.

Remark 3. We refer to [13] on certain correlated random Bernoulli graphs where  $Q^* \neq Q_*$ . On the other hand, there is a class of friendly graphs [1] where  $Q^* = Q_*$ . Existing techniques to improve convex relaxation on GM and QAP include approximate quadratic programming, sorting networks and path following based homotopy methods [12, 18, 21]. Our proposed penalty (1)-(2) appears more direct and generic. A detailed comparison will be left for a future study.

Remark 4. In Example 1, if the convex relaxed  $q_*=2/3$  is rounded up to 1, then  $Q_*=Q^*$ . In Example 2 (Fig. 4), the two interior minimal points at  $\lambda=1.8,1.9$ , after rounding down (up), become zero or one. So convex relaxation with the help of rounding happens to recover the exact permutation. We show in Example 3 below that convex relaxation still fails after rounding (to 1 if the number is above 1/2, to 0 if the number is below 1/2).

**Example 3:** We consider the two-layer neural network model with one hidden layer [11]. Given  $m \ge 0$ , the forward model is the following function:

$$f_m(x, W) = \|\phi((mI + W)x)\|_1,$$

where  $\phi(v) = \max(v, 0)$  is the ReLU activation function,  $x = (x_1, x_2) \in \mathbb{R}^2$  is the input random vector drawn from a probability distribution,  $W \in \mathbb{R}^{2 \times 2}$  is the weight matrix, I is the identity matrix. Consider a two-layer teacher network with  $2 \times 2$  weight matrix

$$W^* = \begin{bmatrix} a & b \\ c & d \end{bmatrix}, \quad a, b, c, d \ge 0.$$

П

We train the student network with doubly stochastic constraint on W using the  $\ell_2$  loss:

$$L(W) = \mathbf{E}_{x} \left[ f_{m}(x, W) - f_{m}(x, W^{*}) \right]^{2}.$$

Let  $p \in [0, 1]$ ,

$$W = \left[ \begin{array}{cc} p & 1-p \\ 1-p & p \end{array} \right].$$

We write the loss function as

$$l_{m}(p) := L(W)$$

$$= E_{x} \left[ \phi \left( (m+p) x_{1} + (1-p) x_{2} \right) + \phi \left( (1-p) x_{1} + (m+p) p x_{2} \right) - \phi \left( a x_{1} + b x_{2} \right) - \phi \left( c x_{1} + d x_{2} \right) \right]^{2}$$

$$= E_{x} \phi \left( (m+p) x_{1} + (1-p) x_{2} \right)^{2} + E_{x} \phi \left( (1-p) x_{1} + (m+p) x_{2} \right)^{2} + 2E_{x} \left[ \phi \left( (m+p) x_{1} + (1-p) x_{2} \right) \cdot \phi \left( (1-p) x_{1} + (m+p) x_{2} \right) \right]$$

$$- 2G_{m} \left( p, a, b \right) - 2G_{m} \left( p, c, d \right) + E_{x} \left[ \phi \left( a x_{1} + b x_{2} \right) + \phi \left( c x_{1} + d x_{2} \right) \right]^{2}, \tag{9}$$

where for  $s, t \ge 0$ ,  $G_m(p, s, t)$  is defined as

$$\mathbf{E}_{x} \left[ \phi \left( (m+p) x_{1} + (1-p) x_{2} \right) \phi \left( s x_{1} + t x_{2} \right) \right. \\ \left. + \phi \left( (1-p) x_{1} + (m+p) x_{2} \right) \phi \left( s x_{1} + t x_{2} \right) \right].$$

Define

$$I(q, r, s, t) := \mathbf{E}_{x} \left[ \phi \left( qx_{1} + rx_{2} \right) \phi \left( sx_{1} + tx_{2} \right) \right],$$

then

$$I\left(q,r,s,t\right)=I\left(s,t,q,r\right),$$

and

$$G_m(p, s, t) = I(m+p, 1-p, s, t) + I(1-p, m+p, s, t).$$

For simplicity, let x obey uniform distribution on  $[-1, 1]^2$ . For

$$qt \ge r$$
,  $q+r > 0$ ,  $s+t > 0$ ,

I(q, r, s, t) equals

$$\begin{cases} \frac{2}{3} (qs+rt) + \frac{q^2 (qt-3rs)}{24r^2} + \frac{s^2 (3qt-rs)}{24t^2}, q < r \\ \frac{1}{3} (qs+rt) + \frac{1}{4} (qt+rs) \\ + \frac{1}{24} (\frac{r^2}{q^2} + \frac{s^2}{t^2}) (3qt-rs), q \ge r \text{ and } t \ge s \\ \frac{2}{3} (qs+rt) + \frac{r^2 (3qt-rs)}{24q^2} + \frac{t^2 (qt-3rs)}{24s^2}, t < s. \end{cases}$$
(10)

We have

$$E_{X}\phi((m+p)x_{1} + (1-p)x_{2})^{2}$$

$$= E_{X}\phi((1-p)x_{1} + (m+p)x_{2})^{2}$$

$$= \frac{2}{3} [(m+p)^{2} + (1-p)^{2}], \qquad (11)$$

$$E_{X} [\phi((m+p)x_{1} + (1-p)x_{2})\phi((1-p)x_{1} + (m+p)x_{2})]$$

$$= \frac{(m+1)^{3}}{3\theta_{m}(p)} + \frac{(m+p)^{4}}{12\theta_{m}(p)^{2}}, \qquad (12)$$

where  $\theta_m(p) := \max(m+p, 1-p)$ . The last term in (9) is a constant:

$$E_{x} \left[ \phi \left( ax_{1} + bx_{2} \right) + \phi \left( cx_{1} + dx_{2} \right) \right]^{2}$$

$$= \frac{2}{3} \left( a^{2} + b^{2} + c^{2} + d^{2} \right) + 2 I \left( a, b, c, d \right). \tag{13}$$

Consider a special case when a = 1/3, b = 2/3, c = 1/4 and d = 3/4. By (11)-(13), the loss function  $l_m(p)$  equals

$$\frac{2}{3}\left[(m+p)^2 + (1-p)^2\right] + \frac{2(m+1)^3}{3\theta_m(p)} + \frac{(m+1)^4}{6\theta_m(p)^2} - 2G_m(p, \frac{1}{3}, \frac{2}{3}) - 2G_m(p, \frac{1}{4}, \frac{3}{4}) + \frac{8113}{5184}.$$

Let

$$F_m(p) := l_m(p) - 4\lambda \sqrt{p^2 + ((1-p)^2)}.$$

When m = 0,  $\lambda = 0$ , Fig. 5 (top) shows  $l_0(p)$  has minimal points in the interior of (0, 1). A permutation matrix W that minimizes L(W) can be achieved by rounding the minimal points. However, when m = 1,  $\lambda = 0$ , (Fig. 5, bottom), rounding the interior minimal point of  $l_1(p)$  gives the wrong permutation matrix at p = 1. At  $\lambda = 0.4$ , the P regularization selects the correct permutation matrix.

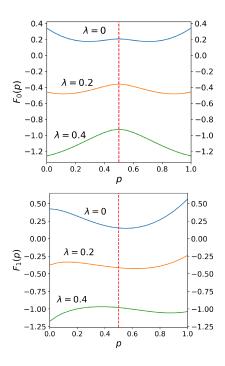


Figure 5:  $F_m(p)$  (m = 0 top, m = 1 bottom) as penalty parameter  $\lambda$  varies for the uniformly distributed input data on  $[-1, 1]^2$ .

Remark 5. If x obeys the unit Gaussian distribution as in [11], the  $F_m(p)$  functions are more complicated analytically, however their plots resemble those for uniformly distributed x, see Fig. 6.

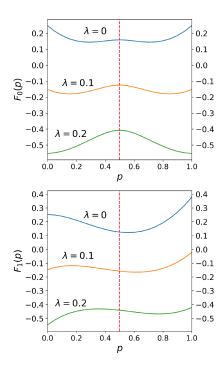


Figure 6:  $F_m(p)$  (m = 0 top, m = 1 bottom) as penalty parameter  $\lambda$  varies for unit Gaussian input data on  $\mathbb{R}^2$ .

# 4 EXPERIMENTS

We relax the shuffle units in ShuffleNet v1 [23] and ShuffleNet v2 [14] and perform experiments on CIFAR-10, CIFAR-100 [9] and a subset of 20 classes (Tab. 3) from ImageNet [5, 10] classification datasets. The 20-class data set consists of most common objects from ImageNet, and forms a typical setting for LCNN application. The accuracy results of auto shuffles are evaluated after the relaxed shuffles are rounded. There is no finetuning of weights after the rounding.

On CIFAR-10 and CIFAR-100 datasets, we set the  $\ell_{1-2}$  penalty parameter  $\lambda=10^{-3}$ . All experiments are randomly initialized with learning rate linearly decaying from 0.2. We train each network for 200 epochs on CIFAR-10 and 300 epochs on CIFAR-100. We set weight-decay  $10^{-4}$ , momentum 0.95 and batch size 128. With w and M initialized from unit Gaussian distribution in Alg. 1, we never run into zero rows and columns in M. An explanation is that those degenerate cases are not generic to cause problems for SGD based training. In Tab. 1-2, we see that auto shuffle consistently improves on manual shuffle in v1 and v2 models of ShuffleNet, by as much as 1.73 % on v1 (g=3). Here g is the number of groups in group convolution. The number of channels is scaled to generate networks of different complexities, marked as 1x, 1.5x, etc.

Next we evaluate auto shuffle in light versions of ShuffleNets (v1 0.25x, v2 0.5x) on a 20-class subset of ImageNet. The subset can be divided into 5 categories, each of which consists of 4 similar classes, see Tab. 3. For each experiment, we set the  $\ell_{1-2}$  penalty parameter  $\lambda=2\times 10^{-3}$ . The training process includes two training cycles: the

Table 1: CIFAR-10 validation accuracies.

Network	v1 (g=8)	v1 (g=3)	v2 1x	v2 1.5x
Manual	90.06	90.55	91.90	92.56
Auto	91.26	91.76	92.81	93.22

Table 2: CIFAR-100 validation accuracies.

Network	v1 (g=8)	v1 (g=3)	v2 1x	v2 1.5x
Manual	69.65	70.16	72.75	
Auto	70.89	71.89	73.40	74.26

first cycle is randomly initialized with learning rate starting at 0.2 and the second one is resumed from the first one with learning rate starting at 0.1. Each cycle consists of 200 epochs and the learning rate decays linearly. We set weight-decay  $4 \times 10^{-5}$ , momentum 0.9 and batch size 128. In Tab. 4, auto shuffle again consistently improves on manual shuffle for both v1 and v2 models, by as much as 2% on v1(g=3). In ShuffleNet v2, r is the fraction of channels that are fed into the right branch of the shuffle unit at Channel Split (Fig. 2). The smaller the r, the lighter the model, the more the auto shuffle improvement.

Table 3: 20-class subset (A=Architectures, L=Landscapes).

Cats	Dogs	Vehicles	A	L
Egyptian	Sheepdog	Bike	Bridge	Valley
Persian	Bulldog	Sports car	Dam	Sandbar
Tiger	Mountain	Scooter	Castle	Cliff
Siamese	Maltese	Cab	Fence	Volcano

Table 4: Validation accuracies on 20-class in Table 3.

Network	v1 (g=8)	v1 (g=3)	v2 (r=0.3)	v2 (r=0.5)
Manual	82.84	82.00	84.63	86.11
Auto	83.68	84.00	85.58	86.84

The permutation matrix of the first shuffle unit in ShuffleNet v1 (g=3) is a matrix of size  $60 \times 60$ , which can be visualized in Fig. 7 (manual, left) along with an auto shuffle (right). The dots (blanks) denote locations of 1's (0's). The auto shuffle looks disordered while the manual shuffle is ordered. However, the inference cost of auto shuffle is same as manual shuffle since the shuffle is fixed and stored after training.

The accuracy drop due to rounding to produce auto shuffle from relaxed shuffle is indicated by relative change in Tab. 5. On CIFAR-10 dataset, negligible drop is observed for ShuffleNet v1. Interestingly, rounding even gained accuracy for on the 20-class dataset.

The  $\ell_{1-2}$  penalty of ShuffleNet v1 (g=3) is plotted in Fig. 8. As the penalty decays, the validation accuracy of **auto shuffle** (after

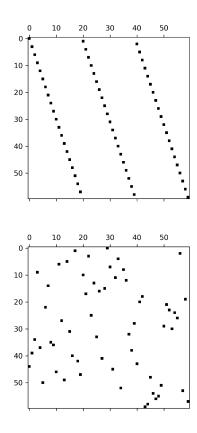


Figure 7: Permutation matrices of the first shuffle unit in ShuffleNet v1 (g=3) of manual shuffle (top) and auto shuffle (bottom). The auto shuffle is trained on CIFAR-10 dataset. The dots (blanks) indicate locations of 1's (0's). The auto shuffle looks disordered while the manual shuffle is ordered. The inference cost of auto shuffle is comparable to manual shuffle in inference.

Table 5: Relative change (Rel. Ch) of accuracy of rounding relaxed shuffle. The -/+ refer to accuracy drop/gain after rounding to produce auto shuffle from relaxed shuffle.

Dataset	Network	Rel. Ch. (%)
	v1 1x (g=8)	0
CIFAR-10	v1 1x (g=3)	0
CIFAR-10	v2 1x	-0.02
	v2 1.5x	-0.11
	v1 0.25x (g=8)	+0.25
20-class	v1 0.25x (g=3)	+0.76
20-Class	v2 0.5x (r=0.3)	+0.37
	v2 0.5x (r=0.5)	0

rounding) becomes closer to **relaxed shuffle** (before rounding), see Fig. 9.

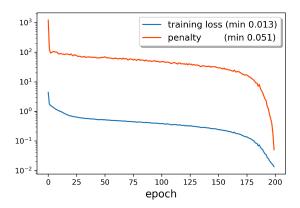


Figure 8: Training loss L and penalty P of ShuffleNet v1 (g=3) with relaxed shuffle on CIFAR-10.

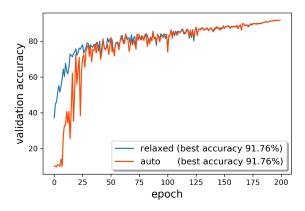


Figure 9: Validation accuracy of ShuffleNet v1 (g=3) with relaxed shuffle (before rounding) and auto shuffle (after rounding) on CIFAR-10. The rounding error becomes smaller during training.

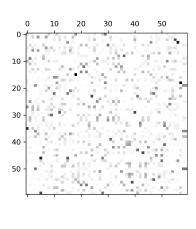
To demonstrate the significance of the  $\ell_{1-2}$  regularization, we also tested auto shuffle with various  $\lambda$  on ShuffleNet v1 (g=3). Tab. 6 shows that the accuracy drops much after the relaxed shuffle is rounded. We plot the stochastic matrix of the first shuffle unit of the network at  $\lambda=0$  and  $\lambda=10^{-5}$  respectively in Fig. 10. The penalty is large when  $\lambda$  is relatively small, indicating that the stochastic matrices learned are not close to optimal permutation matrices.

### 5 CONCLUSION

We introduced a novel, exact and Lipschitz continuous relaxation for permutation and learning channel shuffling in ShuffleNet. The learned shuffle consistently out-performs manual shuffle on CIFAR-10, CIFAR-100, and 20 sub-class of ImageNet data sets across various

Table 6: CIFAR-10 validation accuracies of ShuffleNet v1 (g=3) with relaxed (R) shuffle (before rounding) and auto (A) shuffle (after rounding), and penalty (P) values of relaxed shuffle at various  $\lambda$ 's. The penalty and rounding error tends to zero as  $\lambda$  increases.

λ	0	1E-5	1E-4	5E-4	1E-3
R	90.00	90.18	90.48	91.45	91.76
Α	10.00	38.18	11.37	71.50	91.76
P	3.37E3	1.59E3	4.95E2	3.13E-1	5.07E-2



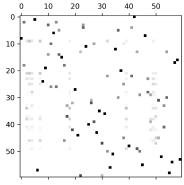


Figure 10: Stochastic matrices of the first shuffle unit in ShuffleNet v1 (g=3) with relaxed shuffle before rounding at  $\lambda=0$  (top) and  $\lambda=10^{-5}$  (bottom). The relaxed shuffle is trained on CIFAR-10 dataset. The matrices are quite diffusive, and not close to optimal permutation matrices when  $\lambda$  is relatively small.

light channel designs while preserving the inference costs of ShuffleNet. We give solvable graph matching examples to show the effectiveness of our permutation penalty. We show analytically through a regression problem of a 2-layer neural network with short cut that convex relaxation of permutation fails even with additional rounding while our relaxation is successful. The idea of auto-shuffle applies broadly to permutation learning problems in science and engineering such as neuron identification from the worm C. elegans [2], image reconstruction from scrambled pieces [4], object tracking [19], to name a few. We plan to extend our work to auto-shuffling in other LCNNs and a wide range of permutation optimization problems of data science in the future.

# **ACKNOWLEDGMENTS**

The work was partially supported by NSF grants IIS-1632935, DMS-1854434, a Qualcomm Faculty Award, and Qualcomm AI Research.

#### REFERENCES

- Y. Aalo, A. Bronstein, and R. Kimmel. 2015. On convex relaxation of graph isomorphism. Proc. National Academy Sci 112(10) (2015), 2942–2947.
- [2] R Badhwar and G Bagler. 2015. Control of Neuronal Network in Caenorhabditis elegans. PLOS One 10(9) (2015).
- [3] R. Burkard. 2013. The quadratic assignment problem. in: Handbook of Combinatorial Optimization (2013), 2741–2814.
- [4] Rodrigo Santa Cruz, Basura Fernando, Anoop Cherian, and Stephen Gould. 2018. Visual Permutation Learning. IEEE Pattern Analysis and Machine Intelligence (2018).
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. 2009. ImageNet: A Large-Scale Hierarchical Image Database. CVPR (2009), 248–255.
- [6] Ernie Esser, Yifei Lou, and Jack Xin. 2013. A Method for Finding Structured Sparse Solutions to Non-negative Least Squares Problems with Applications. SIAM J. Imaging Sciences 6 (2013), 2010–2046.
- [7] Aude Genevay, Gabriel Peyré, and Marco Cuturi. 2018. Learning Generative Models with Sinkhorn Divergences. AISTATS (2018).
- [8] T. Koopmans and M. Beckman. 1957. Assignment problems and the location of economic activities. *The Econometric Society* 25 (1957), 53–76.
- [9] Alex Krizhevsky. 2009. Learning Multiple Layers of Features from Tiny Images. Tech Report (2009).
- [10] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. Imagenet classification with deep convolutional neural networks. NeurIPS (2012), 1097–1105.
- [11] Y. Li and Y. Yuan. 2017. Convergence analysis of two-layer neural networks with ReLU activation. NeurIPS (2017).
- [12] C. Lim and S. Wright. 2016. A Box-Constrained Approach for Hard Permutation Problems. ICML (2016), 10 pages.
- [13] V. Lyzinski, D. Fishkind, M. Fiori, J. Vogelstein, C. Priebe, and G. Sapiro. 2014. Graph matching: Relax at your own risk. arXiv preprint 1405.3133 (2014).
- [14] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. 2018. ShuffleNet V2: Practical Guidelines for Efficient CNN Architecture Design. ECCV (2018).
- [15] G. Mena, D. Belanger, Linderman S, and J. Snoek. 2018. Learning Latent Permutations with Gumbel-Sinkhorn Networks. ICLR (2018).
- [16] Richard Sinkhorn. 1964. A relationship between arbitrary positive matrices and doubly stochastic matrices. The annals of mathematical statistics 35(2) (1964), 876–879.
- [17] K. Sun, M. Li, D. Liu, and J. Wang. 2018. IGCV3: Interleaved Low-Rank Group Convolutions for Efficient Deep Neural Networks. BMVC (2018).
- [18] J. Vogelstein, J. Conroy, V. Lyziński, L. Podrazik, S. Kratzer, E. Harley, D. Fishkind, R. Vogelstein, and C. Priebe. 2015. Fast approximate quadratic programming for graph matching. *PloS one* 10(4) (2015).
- [19] Y. Yang, S. Hallman, D. Ramanan, and C. Fowlkes. 2011. Layered Object Models for Image Segmentation. IEEE Pattern Analysis and Machine Intelligence (2011).
- [20] Penghang Yin, Yifei Lou, Qi He, and Jack Xin. 2015. Minimization of ℓ<sub>1−2</sub> for compressed sensing. SIAM J. Sci. Computing 37(1) (2015), A536–A563.
- [21] M. Zaslavskiy, F. Bach, and J. Vert. 2009. A path following algorithm for the graph matching problem. IEEE Transactions on Pattern Analysis and Machine Intelligence 31 (2009), 2227–2242.
- [22] T. Zhang, G-J Qi, B. Xiao, and J. Wang. 2017. Interleaved group convolutions. CVPR (2017), 4373–4382.
- [23] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. 2017. Shufflenet: An extremely efficient convolutional neural network for mobile devices. CVPR (2017).
- [24] Ritchie Zhao, Yuwei Hu, Jordan Dotzel, Christopher De Sa, and Zhiru Zhang. 2019. Building Efficient Deep Neural Networks with Unitary Group Convolutions. CVPR (2019).